# CamDrop: A New Explanation of Dropout and A Guided Regularization Method for Deep Neural Networks

Hongjun Wang, Guangrun Wang, Guanbin Li, Liang Lin
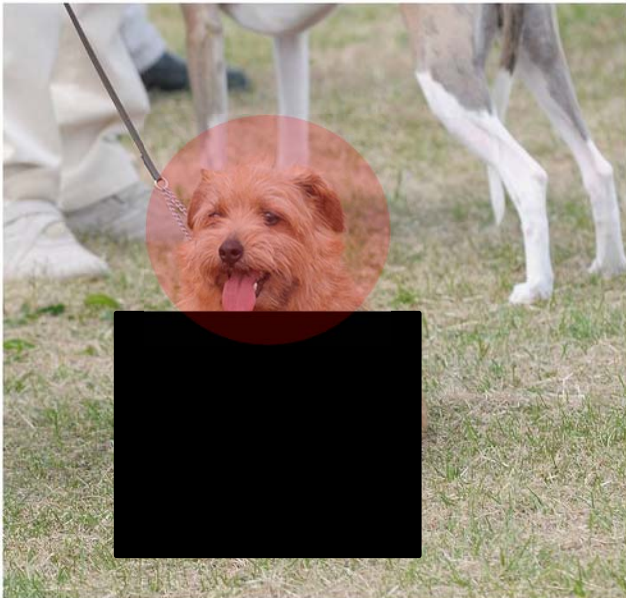
Sun Yat-sen University

- **Introduction**

- Method

- Experiment

- Conclusion

- Take 5 seconds to guess which category these images belong to.

  - It's easy, right?

  - Because you can recognize them by **their distinctive features**

- What about these two?

  – Is it a little hard?

  – But you can still speculate their possible categories by **other neglected parts**
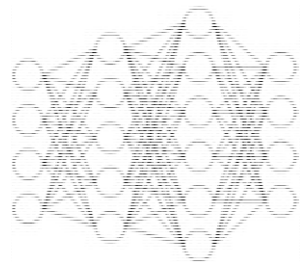


**Cat or Dog**



**Tractor**

- Human race can speculate object categories even though the dominant part is occluded



**How?**

DNN

Human race

- Representations extracted by a robust DNN can **represent more reasonable high-level semantics or detailed spatial information**
- A robust DNN should be **discouraged from "outsmarting" itself by cutting corners**, as clinging to the most obvious point may be one-sided or false positive, which **also known as "overfitting"**
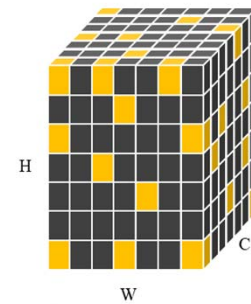
- Existing Dropout methods

  - Recent structured methods are better than others

  - But still lack of semantic guidance

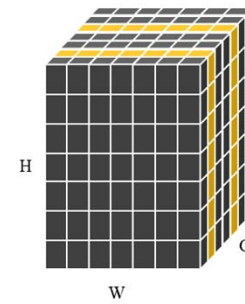$$\widetilde{X}^l = M^l \odot X^l / \gamma$$

$$M^l_{c,h,w} \sim \mathrm{Ber}(\gamma), \qquad \gamma \in [0, 1]$$
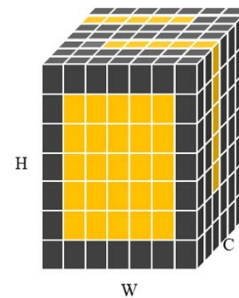
$\gamma$: the retaining rate

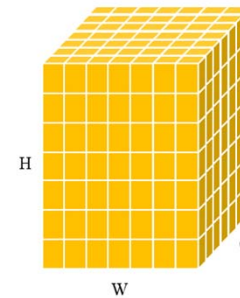$\odot$: the element-wise product of tensors



(a) Dropout  (b) SpatialDropout

(c) DropBlock  (d) DropPath

- Main idea

  - Utilize Class Activation Map (CAM) as the semantic guidance to *hides discriminative object parts* during training

  - Force the network to proactively *explore other neglected parts autonomously* instead of relying on external data

- Introduction

- **Method**
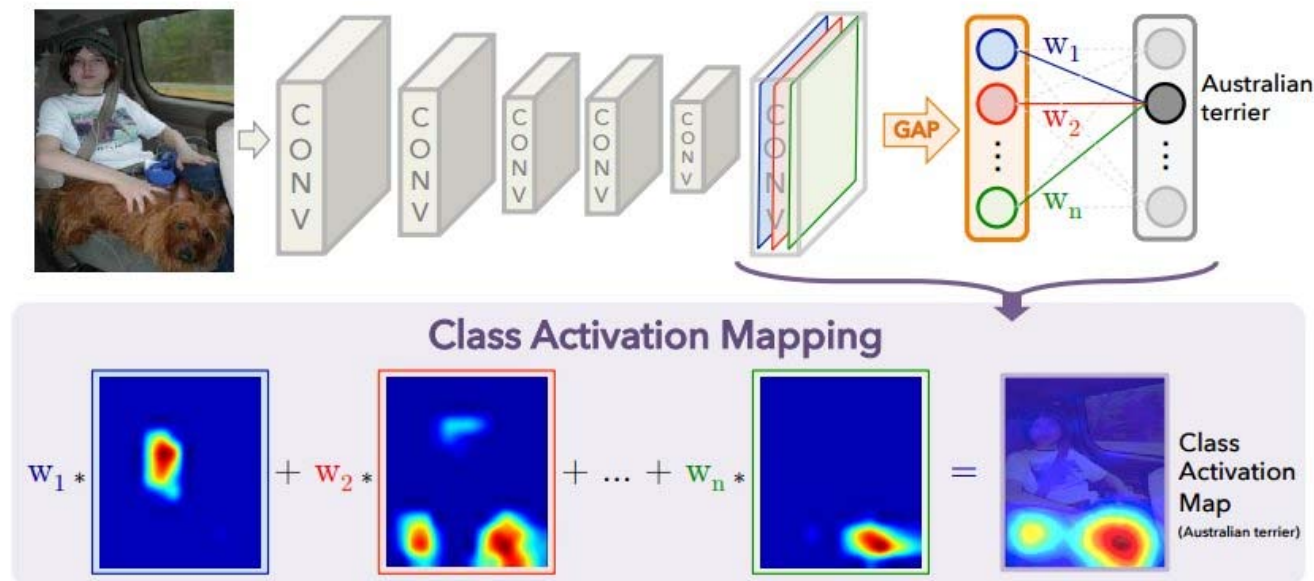
- Experiment

- Conclusion

- Class Activation Map (CAM)
  - You can see that the CNN is triggered by different semantic regions of the image
  - CAM indicates **the discriminative image regions used by the CNN to identify the particular category**



Class Activation Mapping

B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization.

# (0) Generating **class activation map** $J^{k'} \in \mathbb{R}^{1 \times H' \times W'}$:

### i. GAP layer:

$$J^{k'} = \sum_{c=1}^{C'} \alpha_c^{k'} Z_{c,:,:}$$

### ii. Shallow Layers:

$$J^{k'} = \sum_{c=1}^{C''} (\sum_{j=1}^{\lfloor \frac{C'}{C''} \rfloor} \alpha_j^{k'}) X_{c,:,:}^{l}$$

where the weight vector $\alpha^{k'} \in \mathbb{R}^{C' \times 1}$ indicates the importance of each channel for a specific class $k'$ of y, Z, X are the output feature map of the penultimate/middle layer.



(a) Dropout  (b) SpatialDropout

(c) DropBlock  (d) DropPath

(e) **CamDrop**

Importance Weights

Class Activation Mapping

## (1) Spatial-wise binary mask $M^{(1)}$:

$$M^{(1)}_{h,w} = \begin{cases} 0, & J^{k'}_{h,w} > \inf\{T^s_n\} \\ 1, & J^{k'}_{h,w} < \inf\{T^s_n\} \end{cases}$$

$J^{k'}$ can be considered as a set with $H' \times W'$ elements $\{j_{1,1}, \dots, j_{H' \times W'}\}$, each of which implies the significance of units at spatial grid $(h, w)$.

$T^s_n$ is the set of **the $n$ most important pixels**.



(a) Dropout    (b) SpatialDropout

(c) DropBlock    (d) DropPath

(e) **CamDrop**

Importance Weights

Class Activation Mapping

(2) Depth-wise binary mask $M^{(2)}$:

$$M_c^{(2)} = \begin{cases} 0, & \alpha_c^{k'} > \inf\left\{T_{n'}^d\right\} \\ 1, & \alpha_c^{k'} < \inf\left\{T_{n'}^d\right\} \end{cases}$$

where $n'$ is a hyperparameter and $\inf\{\cdot\}$ is the infimum of a set
$M_c^{(2)}$ is the set of **the $n$ most important channels in $\alpha^{k'}$**



(a) Dropout (b) SpatialDropout (c) DropBlock (d) DropPath (e) **CamDrop**

Importance Weights

Class Activation Mapping

(3) Let $M^{(3)}$ be the valid seed region of the feature map and $\psi$ be a randomly sampled mask.
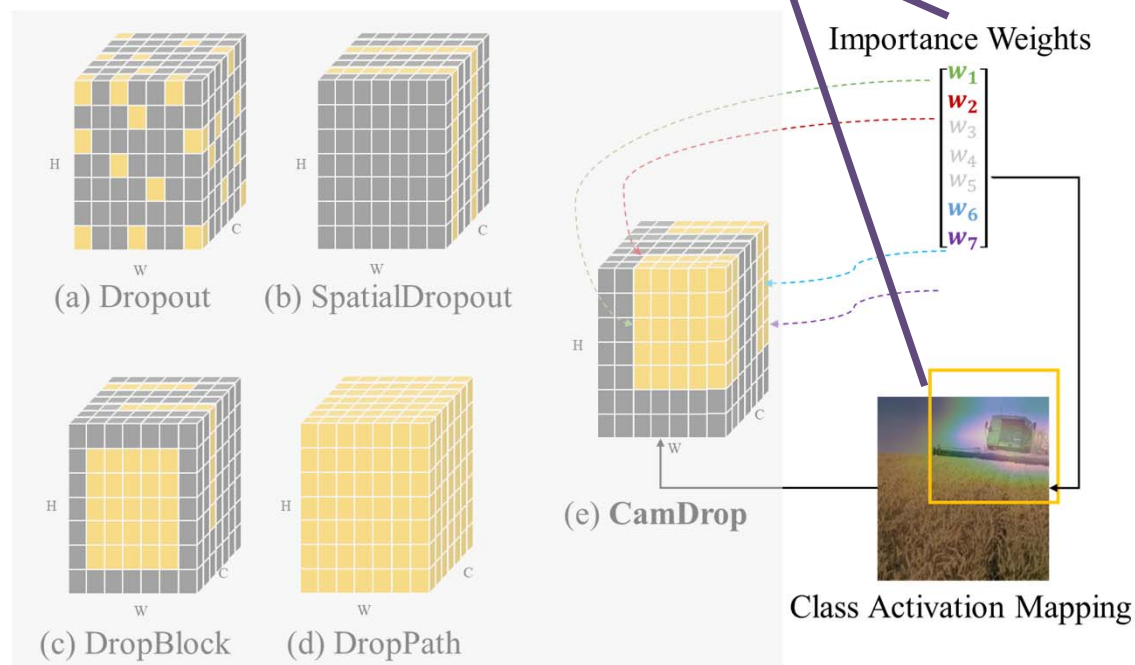
The combined mask $\mathcal{M}$ is:

$$\mathcal{M} = \neg\boldsymbol{\psi} \wedge \left(\bigwedge_i \neg M^{(i)}\right)$$

$\psi \sim Bernouli(\gamma)$



Importance Weights

$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \\ w_7 \end{bmatrix}$

H — C — W
(a) Dropout

H — C — W
(b) SpatialDropout

H — C — W
(c) DropBlock

H — C — W
(d) DropPath

H — C — W
(e) **CamDrop**

Class Activation Mapping
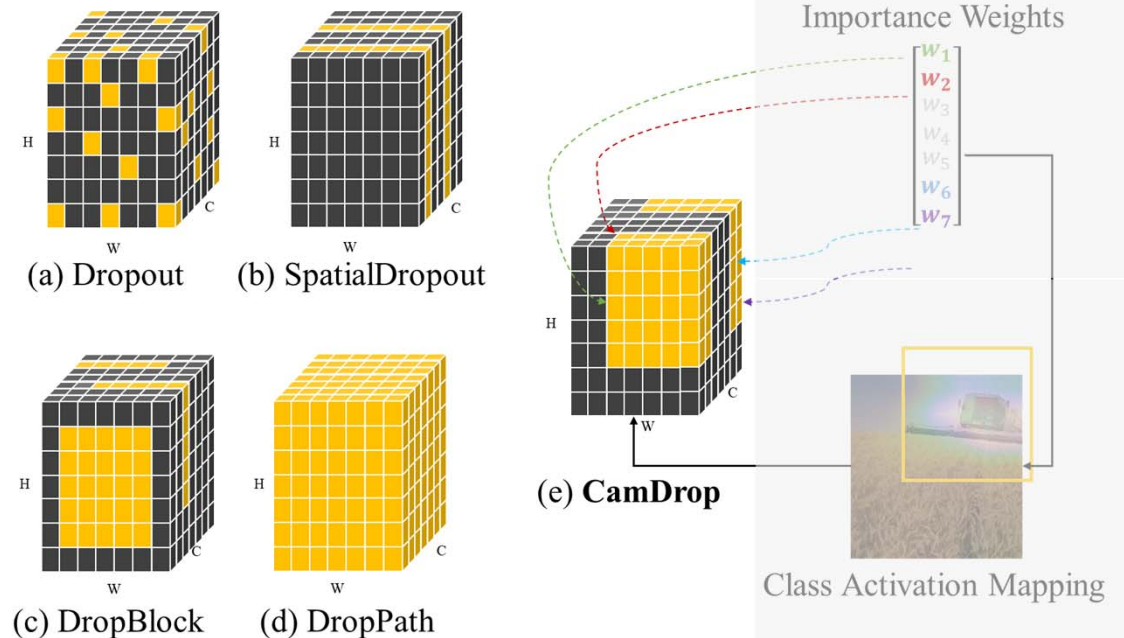
(4) Sweep over $\mathcal{M}$ and set pixels in $\{\forall q, \|q - u\|_1 = r\}$ zero, where $u$ is the zero entries in $\mathcal{M}$ and $r$ is the length of mask.

We normalize $\mathcal{M}$ by the factor:

$$\mathbf{M}^l = \frac{CHW}{\sum_{c=1,h=1,w=1}^{C,H,W} \mathcal{M}_{c,h,w}}$$



(a) Dropout   (b) SpatialDropout

(c) DropBlock   (d) DropPath

(e) **CamDrop**

Importance Weights

Class Activation Mapping

- Introduction
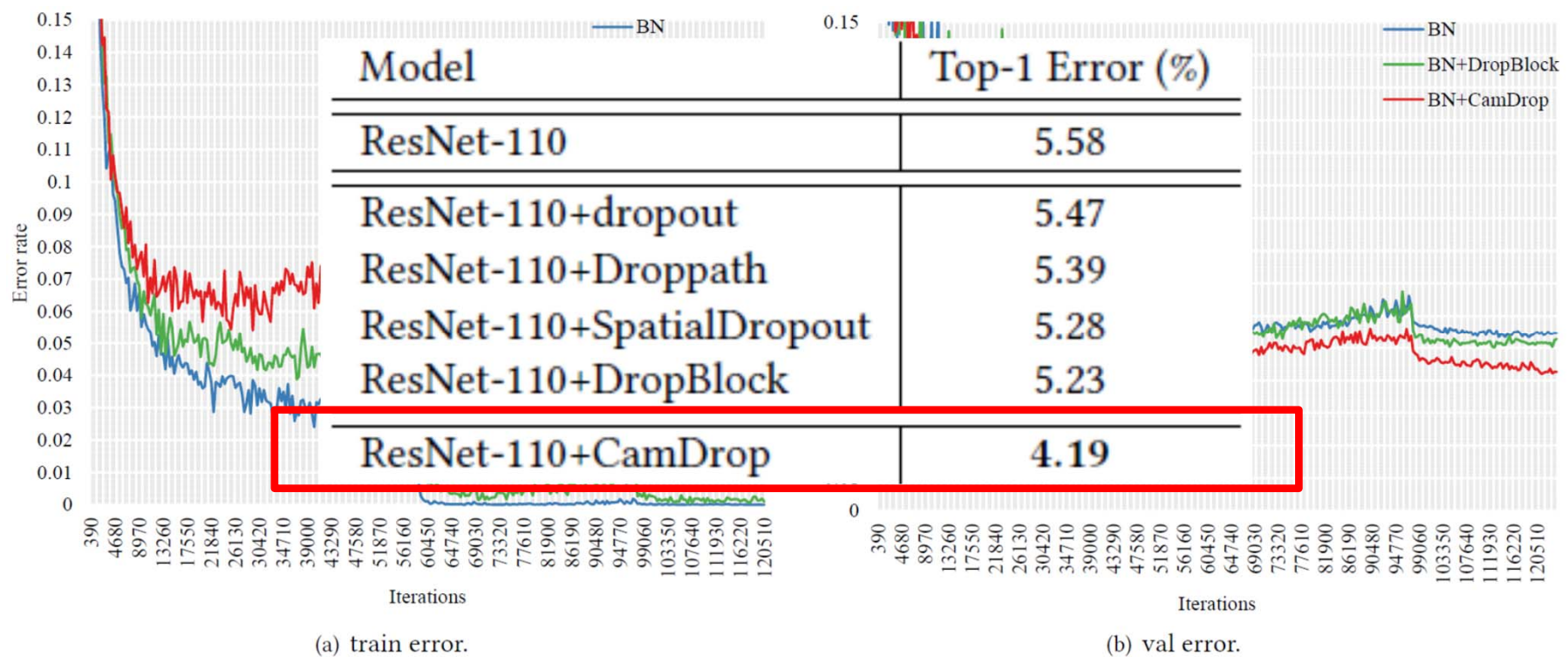
- Method

- **Experiment**

- Conclusion

- Main Experiments
  - CIFAR10
  - ImageNet2012

- Ablation Studies
  - Effectiveness of CAM
  - Effectiveness of Importance Weights
  - Proportion of Dominant Visual Patterns

- More than Overfitting

- Classification on CIFAR10



| Model | Top-1 Error (%) |
|---|---|
| ResNet-110 | 5.58 |
| ResNet-110+dropout | 5.47 |
| ResNet-110+Droppath | 5.39 |
| ResNet-110+SpatialDropout | 5.28 |
| ResNet-110+DropBlock | 5.23 |
| ResNet-110+CamDrop | 4.19 |

(a) train error.

(b) val error.

- Classification on ImageNet



(a) Original  (b) DropBlock  (c) CamDrop

- ## Effectiveness of CAM

Average of feature map



Importance Weights

(a) Dropout  (b) SpatialDropout

(c) DropBlock  (d) DropPath

(e) **CamDrop**

Class Activation Mapping

| | Top-1 Error (%) |
|---|---|
| Avg guided | 5.07 |
| CAM guided | **4.19** |

- Effectiveness of Importance Weights



(a) Dropout   (b) SpatialDropout

(c) DropBlock   (d) DropPath

(e) **CamDrop**

Importance Weights

$w_1$
$w_3$
$w_4$
$w_5$
$w_6$
$w_7$

Class Activation Mapping

| | Top-1 Error (%) |
| --- | --- |
| All | 4.51 |
| Dominants | **4.19** |

- Proportion of Dominant Visual Patterns



(a) Dropout  (b) SpatialDropout

(c) DropBlock  (d) DropPath

(e) **CamDrop**

Importance Weights

Class Activation Mapping

| | Top-1 Error (%) |
|---|---|
| $n' = C/16$ | 4.39 |
| $n' = C/8$ | **4.19** |
| $n' = C/4$ | 4.75 |
| $n' = C/2$ | 4.78 |

- More than dealing with Overfitting
  - When inverting the decay scheme of learning rate, the validation errors of models trained with CamDrop are still close **even though the inversed one is overfitting** (**2nd/3rd rows**).

  - **The two overfitting with and without CamDrop** have a large gap (**1st/3rd rows**).

| Model | Val Error (%) | Train Error (%) |
|---|---|---|
| ResNet-110, $lr = 2.0$ | 7.16 | 2.22e-4 |
| ResNet-110, $lr = 2.0$ with CamDrop | 5.94 | 8.53e-2 |
| ResNet-110, $lr = 2.0$ with CamDrop, inversely | 6.29 | 1.66e-4 |

At the $t$-th update iteration in SGD, the weights and biases in the $l$-th layer will be updated as:

$$\mathbf{b}_t^l := \mathbf{b}_{t-1}^l - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}_{t-1}^l}$$

According to the chain rules:

$$\frac{\partial \mathcal{L}}{\partial b_{t-1,i}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{S}} \frac{\partial \mathbf{S}}{\partial b_{t-1,i}^l}$$

The upper bounds of the loss can be given by the Hölder inequality:

$$\left| \frac{\partial \mathcal{L}}{\partial b_{t-1,i}^l} \right| \leq \max_{k'} \left| \frac{\partial \mathcal{L}}{\partial S_{k'}} \right| \left\| \frac{\partial \mathbf{S}}{\partial b_{t-1,i}^l} \right\|_1$$

Since $|\partial L/\partial S| = |softmax(S) - y|$ cannot exceed 1 for any element, the inequality can be reduced to:

$$\left| \frac{\partial \mathcal{L}}{\partial b_{t-1,i}} \right| \leq \left\| \frac{\partial \mathbf{S}}{\partial b_{t-1,i}} \right\|_1$$

CamDrop masks out the several notable neurons, which gives a **tighter upper bound** of the update of biases $\partial L/\partial b$ at each iteration:
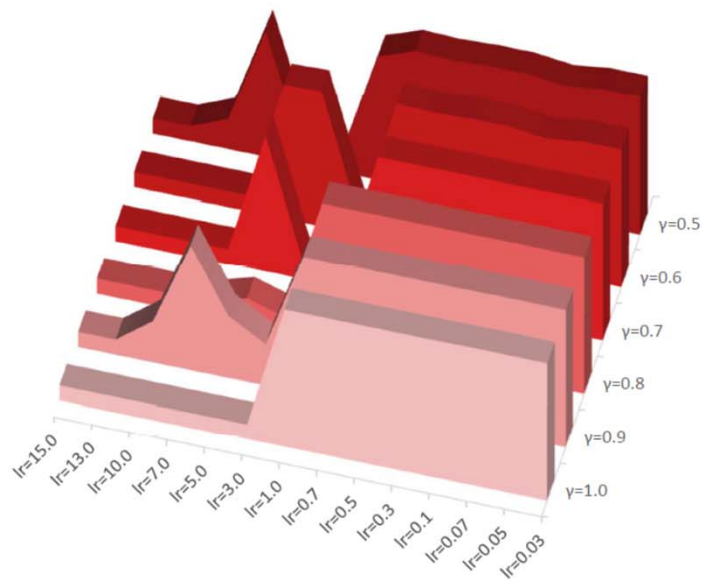
$$\left| \frac{\partial \mathcal{L}}{\partial b_{t-1,i}^l} \right| \leq \left\| \frac{\partial \mathbf{S}}{\partial b_{t-1,i}^l} \right\|_1^{(c)} \leq \left\| \frac{\partial \mathbf{S}}{\partial b_{t-1,i}^l} \right\|_1^{(t)} \leq \left\| \frac{\partial \mathbf{S}}{\partial b_{t-1,i}^l} \right\|_1$$

where $\|\cdot\|^{(c)}$ and $\|\cdot\|^{(t)}$ stand for the $L_1$-norm of derivative with Cam/traditional dropout mask respectively
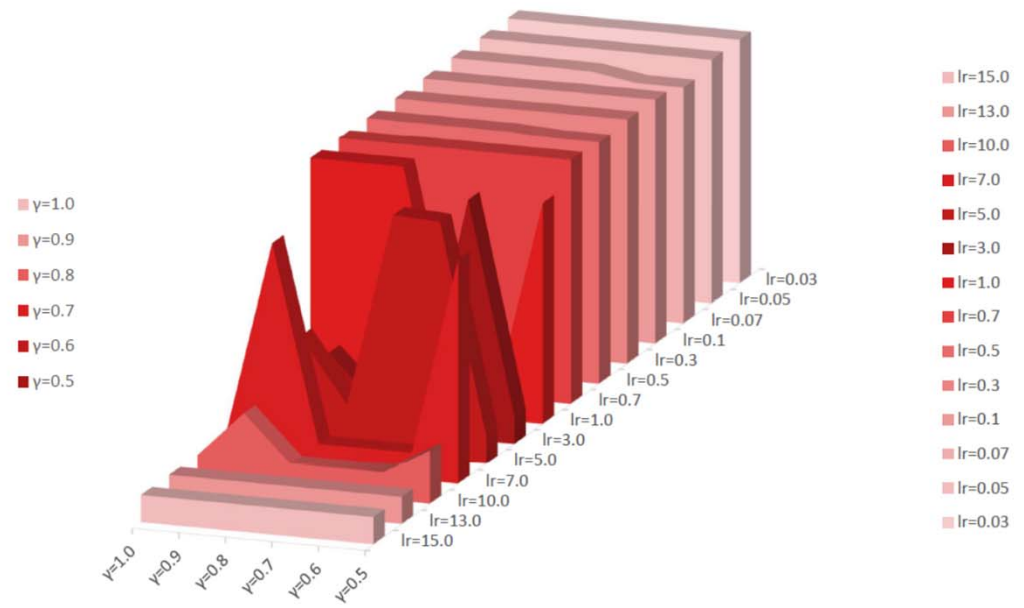
- Relationships among the learning rate, the remaining rate and top1-accuracy (z-axis)



(a)

(b)

- Introduction

- Method

- Experiment

- **Conclusion**

- Conclusion
  - A novel dropout method called **CamDrop** to improve the robustness of DNN models by **masking the dominant regions** with the guidance of **class activation mapping**.

  - Dropout techniques actually make **the upper bound of the magnitude of gradients much tighter**.

  - Data with **non-Euclidean structure** can utilize this technique by establishing graph with corresponding relationships between vertices and edges.

# Thank you

**Mail: wanghq8@mail2.sysu.edu.cn**